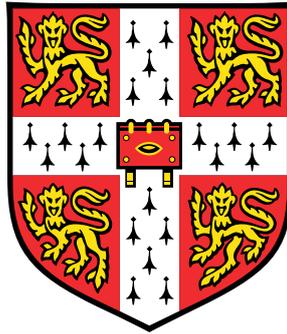# Genomic diversity and speciation in East African cichlid fish

**Milan Malinsky**

Wellcome Trust Sanger Institute

University of Cambridge

This dissertation is submitted for the degree of

*Doctor of Philosophy*

Girton College                    September 2015

To Alena and Sasha . . .

# Declaration

I hereby declare that except where specific reference is made to the work of others, the contents of this dissertation are original and have not been submitted in whole or in part for consideration for any other degree or qualification in this, or any other university. This dissertation is the result of my own work and includes nothing which is the outcome of work done in collaboration, except where specifically indicated in the text. This thesis does not exceed the length limit of 60,000 words specified by the Biology Degree Committee.

Milan Malinsky
September 2015

# Acknowledgements

I would like to acknowledge the people who have helped along the way to becoming a scientist. First my parents and especially my grandparents who encouraged my natural curiosity from an early age. They did not burden me by making me believe in given wisdoms, but instead led me to follow my observations and analysis and make my own conclusions about the world we live in. For this, I am and always will be grateful.

I would like to thank my teachers in the Czech Republic, the most important among them Martina Aberlová whose influence in the final three years of secondary school helped me to get a well-rounded basic education, and at high school Alena Šteklová who encouraged me to continue learning the English language. This enabled me to enter the English speaking world and thus the global scientific enterprise.

Many specific skills I used during my PhD, especially in computer programming and algorithms, were developed during an undergraduate degree in Computer Science at the University of Birmingham, UK. Here, my thanks go to William Edmondson, the Admission Tutor for his confidence in me, Peter Coxhead, my personal tutor, for his help and guidance, and Hamid Dehghani, my final year project supervisor, for introducing me to the world of cutting edge science at an international conference.

Next, my thanks go to Stephen Eglen and Simon Tavaré, for running the Cambridge MPhil course in Computational Biology and for helping me to make the most of it. It was while I was a student on this course that Simon introduced me to Eric Miska who, fascinated by the extraordinary evolution and diversity of East African cichlids, provided the initial spark for this PhD project.

Soon, I too became excited by cichlid evolution and could see that computational analysis of genomic data would provide important clues about the amazing history of this fish family. My PhD supervisor Richard Durbin is a research leader in analysing population genomic data and his help has been invaluable to this project. Access to the Sanger Institute facilities and Richard's expert strategic guidance are what made this project truly possible. I am also very appreciative of the specific insights that Richard provided during the project, often helping to push forward when working on my own felt like I was reaching a dead-end. In the first year in Richard's group I especially

enjoyed working with Jared Simpson, who introduced me to algorithms for genome assembly and to C++ programming. In the next three years, I enjoyed collaborating with Stephan Schiffels with whom I had numerous discussions of population genetics topics, each of which increased my understanding of this scientific field.

George Turner and Martin Genner have been my principal collaborators and advisors on evolution and ecology of cichlid fish. They helped me to develop from being a 'Computer Scientist' at the beginning of the PhD to the 'computer-savvy' evolutionary biologist I am today. Without their input this work and my personal development would be greatly diminished.

I owe thanks to the Wellcome Trust for financially supporting my PhD.

# Abstract

Unravelling the genetic basis of functional diversification is fundamental for our understanding of vertebrate evolution and can also have significant implications for animal and human health. Speciation leads to phenotypic diversity by producing new units of evolution - species. In less than five million years, East African cichlids have radiated into thousands of species that differ in craniofacial morphology, pigmentation, behaviour and many other traits. In this thesis, I take advantage of recent advances in DNA sequencing technologies to study the genetic basis of this exceptional diversity. First, as a member of the Cichlid Genome Consortium (CGC), I identified and characterised over 1,000 loci generating microRNAs, non-coding RNA genes that regulate expression and may play a role in the evolution of cichlid traits. Next, at the Sanger Institute, we obtained whole genome sequences of 271 individuals from over 70 species from in and around Lake Malawi. I aligned the data to a reference genome generated by CGC, and used the results to: 1) ascertain the overall levels of genetic variation and allele sharing within and between species; 2) reconstruct relationships between the species; 3) study in detail the genetic causes and consequences of early stages of speciation in Lake Massoko, a small isolated crater lake in southern Tanzania. I found that that the genetic distance between the most diverged Lake Malawi species is surprisingly low, comparable to the distance between two strains of zebrafish, that there are discrepancies between relationships inferred from molecular phylogeny and from traditional taxonomy, and that measurable introgression between species occurs but does not seem to be common. In Lake Massoko, I identified clearly demarcated genomic regions of differentiation between incipient species in sympatry. Interestingly, there are no fixed differences; instead I found a genome-wide pattern with dozens of loci of moderate divergence. With collaborators, we found that alleles in the regions are associated to mate preferences in the laboratory, and genes in the regions are enriched for molecular functions consistent with morphological and sensory system adaptation. To facilitate this work, I constructed whole genome alignments between CGC genome assemblies, assigned ancestral alleles to genetic variants in Lake Malawi, and built a genome browser that can be used to visualise datasets produced by us and

the CGC. The browser website has been visited over 650 times since March 2014. In addition, I developed a new method for genome assembly to reduce problems caused by heterozygosity, taking advantage of mother-father-offspring trio data. I applied this method to obtain *de novo* genome assemblies of three cichlid species, and also three highly heterozygous *Heliconius* butterfly species. These datasets, tools, and findings make significant contributions to evolutionary genetics and will provide a foundation for future research on processes underlying the evolution of phenotypic diversity, especially in cichlids.

# Table of contents

# List of figures

# List of tables